

POSTED PRECHARGE AND MULTIPLE OPEN-PAGE RAM ARCHITECTURE

DESCRIPTION

[Para 1] CROSS REFERENCE TO RELATED APPLICATIONS

[Para 2] This application claims the benefit of U.S. Provisional Application Nos. 60/510,182 filed October 10, 2003, and 60/525,453, filed November 28, 2003.

[Para 3] BACKGROUND OF THE INVENTION

[Para 4] [0001] The present invention generally relates to electronic systems utilizing memory devices. More particularly, this invention relates to a method and memory architecture that enables multiple pages to be open on the same bank of a memory device.

[Para 5] [0002] High density random access memories (RAM) find numerous applications in computers, communications, and consumer and industrial applications. These memories are mostly volatile (lose stored data when power is switched off). These cost sensitive, medium performance (high bandwidth with moderate latency) segments have been served by multibank SDRAM's (synchronous dynamic random access memories). SDR SDRAM's, DDR-I SDRAM's, DDR-II SDRAM's, QDR SDRAM's, Rambus™ DRAM's, Network DRAM's™, FCRAM's™, and RLDRAM's™ are among the numerous varieties of DRAM available in the commercial market. In addition to DRAM, other types of memory finding use in a variety of applications include high-density flash

memory products, which provide nonvolatile storage at reasonably high densities. High-density flash memories are used in portable electronic appliances such as cell phones, digital cameras, etc.

[Para 6] [0003] Emerging applications in computer and communications processors (CPU's, NPU's) will increasingly require (in addition to bandwidth and fairly low latency at low cost) the ability to open new rows in the currently active banks in these SDRAM's (single-chip IC's, DIMM modules, RIMM modules, or other memory subsystems comprising multiple SDRAM's) for superior performance. Specifically, the advent of hyperthreading, multithreading, multichannel-DRAM access, shared memory and other processor inventions is making essential a different SDRAM architecture with a multiple open-page capability. However, current DRAM architectures and timings require "closing current active row/page in a bank," before a new row/page can be "opened" in the same bank.

[Para 7] [0004] Current leading-edge general-purpose SDRAM's are multibank devices designated as "DDR-I" or just "DDR" (double data rate) for the current generation of devices and "DDR-II" or "DDR2" for the emerging new generation. Four internal banks are the current standard; eight-bank DDR-II devices are now being introduced. Herein, "DDR" refers generically to both DDR-I and DDR-II, unless otherwise indicated. References to SDRAM's are to all SDRAM's including but not limited to DDRs, unless otherwise indicated.

[Para 8] [0005] As is true of SDRAM's in general, and assuming proper memory controller support, any one or more of the internal banks can be activated by an Active command at one (but only one) row address at a given time. The row addresses can be the same or different for any or all of the banks. This basic "one open row per bank at a time" limitation arises from the inherent properties of the one-transistor, one-capacitor (1T/1C) storage cells that comprise the basic storage unit of the SDRAM core, combined with the

peripheral architecture of the multibank SDRAM. If a memory subsystem has memory devices mounted in separate “ranks” (on a module) or in two (or possibly more) modules inserted into separate slots on the motherboard, the capacity of the memory is correspondingly increased and the number of open pages can be increased to a maximum of the number of ranks times the number of internal banks per device. Under that peripheral architecture, all banks have to share the same, single data path to the external world. Even with separate data I/O circuitry for each bank, it is still only one bank at a time that can accept Read or Write commands from the memory controller. The 1T/1C memory cell, which is the essence of SDRAM architecture in general for cost and density reasons, imposes the further “one open row per bank” limitation. These limitations apply to SDR SDRAM’s as well as DDR SDRAM’s, and other DRAM’s using 1T/1C memory cell as the basic storage element. In addition to “restore after read,” a sense/read/write amplifier dedicated to precharge of bit line(s) as well as sense/restore functions aggravates the latency limitation.

[Para 9] [0006] In order to activate a bank for reading from or writing to one of its rows, the bank must first be charged in order to allow the chip to sense the particular row and amplify the signal from that row. Rather than waiting for a Read/Write request before being charged, a memory bank is usually precharged, initiated by a “Precharge” command. Opening a row in an SDR or DDR SDRAM activates a “page,” which is a 2-dimensional array of bits defined by the length of the row in one dimension, and the “width” of the device (being equal to the number of I/O pins in the device) in the second dimension. The length of the row is determined by the number of columns that intersect the row. In total number of bits, the size of the page is, therefore, equal to the number of column addresses multiplied by the number of I/O pins. Whenever data is read from or written to a device, the specific column number within a row address is furnished and a string of bits equal to the number of I/O pins can be retrieved from or stored in the device. The number of column addresses equals 2 to the power of the number of column

address lines. This I/O is combined with I/O from or to the other devices on the memory module(s) comprising the SDRAM memory subsystem, so as to deliver the total package of bits (generally 64 bits in the case of conventional SDRAM memory subsystems, though with important exceptions) onto the memory bus in the case of a Read or store it in the module(s) in the case of a Write.

[Para 10] [0007] Closing of any bank is initiated by a Precharge command, which means that the command bus will be busy during the Precharge command. In SDRAM's in general, precharging introduces a delay or latency before the new bank can be given an Activate command whereby the process of opening the new bank is started. This "Precharge Latency" is often referred to symbolically as tRP, and is typically in the range of 2 to 5 clock cycles. There are numerous timing parameters that define the actual performance of SDRAM's in general and DDR-II devices in particular. The interaction of these parameters varies considerably, sometimes greatly, depending on the exact circumstances of a given access to the memory subsystem. The interaction is extremely complex, and requires many timing diagrams for a rigorous explanation. Only certain timing parameters will be discussed herein, and their explanations will be necessarily simplified and instructive only in its coverage of some of the various combinations of the timing parameters. Even if the Precharge can be hidden behind an ongoing output burst, the command bus shared by all devices will be busy, and thus prevent the issuing of other commands. Conflicting needs for commands during the same cycle are usually referred to as "command bus contention" and are a significant component of tRP and other limitations.

[Para 11] [0008] Following satisfaction of this latency and upon issuance of the Activate command, there begins a second latency period, often described as the "Active to Read or Write command delay," or "RAS-to-CAS delay," or a similar term, and is generally symbolized by tRCD. This latency is generally in the range of 2 to 5 clocks, and must be satisfied before a Read or

Write command can be issued. In order to avoid command bus contention, starting with DDR-II, the Read or Write command can be issued as early as the next clock cycle following the Activate command. However, since tRCD needs to be satisfied before the execution of the Read or Write command can commence, the command needs to be pushed out internally within the device according to a predefined “Additive Latency” (AL). This mode of operation is referred to as postponed or “posted” CAS mode.

[Para 12] [0009] Following this comes another latency period, generally known as “/CAS Latency” (CL), a programmed value generally 2 to 4 clocks in length, which begins with the commencement of internal execution of the Read or Write command, and extends until the data is placed on the memory bus for sending to the processor in the case of a Read or written into the SDRAM in the case of a Write. In the nomenclature of the current JEDEC DDR-II Standard, and commonly in industry specification sheets also, the sum of AL and CL equals the “Read Latency” or “Read Access of the first Critical Word.” The equivalent latency parameter for Writes, “Write Latency,” equals Read Latency minus 1.

[Para 13] [0010] In a DDR device, depending on speed grade and other factors, the access latencies typically total in the range of 4 to 7 or more clocks to which the Precharge latencies (2 to 4 clocks) need to be added. The total is generally referred to as the “Active-to-Active command interval – Auto-Precharge,” the “bank cycle time,” or a similar term, generally symbolized as tRC. In short, the bank cycle time refers to the number of clocks required between two consecutive accesses to different rows in the same Bank, irrespective of whether they are Reads or Writes.

[Para 14] [0011] To be 100% efficient, an SDRAM memory subsystem would need to deliver data to the memory bus (in the case of Reads) or to the SDRAM array (in the case of Writes) on every clock cycle so as to respond fully

to the processor's need to load (Read) or store (Write) data. In the case of DDR, the requirement is to deliver 2 bits of data per clock cycle at each I/O pin in the memory subsystem. If it took 9 to 12 or more clocks to deliver this quantity of data, efficiency would be near or below an unacceptable 10%. To improve on this in the case of consecutive accesses to the same bank, the bank is left activated/open after the first access, without Precharge. The first access to the bank incurs the full latencies, but accesses to the same bank thereafter can be made contiguous or nearly so by using the "burst" technique. Under this technique, just the starting column address for a burst of 2, 4 or 8 column addresses is strobed (4 or 8 for the DDR-II). The tRP and tRCD latencies are not incurred again as long as the row/page remains open, and because of the extremely high speed of internal SDRAM operation, successive column accesses equal to the length of the burst are enabled. Since only a single column address is strobed, /CAS Latency also need not be incurred again. Using this technique, once all latencies have been satisfied one time and accesses thereafter are confined to sequential column segments in the same row/page, the device can keep pace with the processor's demands for Reads and Writes and processor clock cycles are not wasted. If the device is programmed for "interleaved" burst accesses, the columns within a burst sequence are actually accessed in a numerically non-consecutive, although pre-determined, order. If the burst length is 4 or 8, as in the case of a DDR-II device, subsequent commands and their related latencies (including tRP when a Precharge command is involved) can be partially or (especially if the burst length is 8) entirely hidden behind the burst operation.

[Para 15] [0012] On the basis of the above, largely uninterrupted back-to-back transactions to or from the processor are possible, assuming the bursted data is fully usable by the processor. However, this assumption is not correct in many cases. Typically, a data word is 8 Bytes (64 bits) wide from a memory DIMM or RIMM. Having 8-bit bursts necessitates instruction execution on 64 Bytes of data in the SAME PAGE to take maximum advantage – a very unlikely event in most applications. Moreover, in the case of burst Writes, a Write

Recovery Latency (tWR) typically of 3 clocks after the last burst Write access must be satisfied before there can be issued a subsequent Precharge command, which must precede the access of any other bank. This is a “non-hideable” 3-clock latency imposing a 6-bit interruption (in the case of DDR) per I/O pin of the data flow, a significant performance loss. Before the current leading edge DDR-II, the AC operating characteristics of the DDR SDRAM were somewhat simpler and entailed use of several fewer timing parameters, but the basic issues described above also arose in substantially the same way.

[Para 16] [0013] If the row/page to be newly opened is located in a different bank of the device, the time period between the Activate commands in the old and new banks is often called the “Active bank A to active bank B command period” or “Row-to-Row Delay” (tRRD), which is needed to satisfy the row decoder latency. tRRD specifies the latency of jumping from one open row to another open row of different banks within the memory subsystem space. tRRD is typically on the order of 2 clocks and can usually be hidden behind ongoing data bursts, and therefore is likely to constitute a problem in cases of random, single-word accesses (but not otherwise).

[Para 17] [0014] Current SDRAM’s do offer a programmable Read or Write Auto-Precharge function which in some cases reduces the tRP latency by hiding part of it. However, this is at the cost of closing the currently open row, thereby negating the advantage of the open page architecture, namely the extremely fast Reads/Writes of sequences of data which are located contiguously (i.e., in the same open page) in the SDRAM. Therefore, while the Auto-Precharge function can hide tRP latency and permit faster random accesses, it also defeats the principal advantage of the open page organization.

[Para 18] [0015] In view of the above, the performance of RAM devices can be limited by the inability to hold multiple pages open on the same bank of the

memory subsystem. A worst case scenario would be encountered if data from two pages are requested where, after closing the first page and moving to the second page, the first page is needed again for additional data. The possibility of incurring this access pattern also constitutes the main performance limitation of pseudo-SRAM (pseudo static random access memory, or PSRAM) since it disallows single random accesses without satisfying the minimum RAS pulse width. Flash memory devices, both NOR and NAND type, also use a “page” architecture, for accessing the memory device. Therefore, as with DRAM devices, flash memory devices also suffer from the inability to hold multiple pages open on a single bank.

[Para 19] BRIEF SUMMARY OF THE INVENTION

[Para 20] [0016] The present invention provides a method and architecture that overcomes the problem of latency-caused performance degradation discussed above. The invention involves the use of what is termed herein a “Posted Precharge,” which as used herein means that an external command for Precharge is given as early as possible, preferably immediately following a Read command. However, the present invention causes the execution of the Precharge to be delayed by means of an internal mechanism, preferably a counter means, until all Read/Write commands are completed. According to a preferred aspect of the invention, the method and architecture also provide that any subsequent access to the same page causes the counter means to be reset, further delaying the execution of the Precharge operation. In other words, the invention provides an internal Precharge scheme that starts the Precharge for that particular bank, activity-dependent, at an appropriate time, thus freeing up the command bus.

[Para 21] [0017] In view of the above, the Posted Precharge of the present invention finds analogy in some respects to the posted Read/Write command architecture of DDR-II SDRAM’s referred to above. Moreover, compared to the

simpler current scheme of an Auto-Precharge following a Read or Write, the Posted Precharge of this invention has the flexibility to accommodate single random accesses as well as bursts. Notably, the DDR-II architecture does not allow row accesses of less than a burst length of 4.

[Para 22] [0018] The benefits of this invention can be promoted by implementing the Posted Precharge on the DDR device itself, which has the potential for virtually eliminating command bus contentions. The benefits of Posted Precharge can also be obtained in DIMM's and SIMM's (and similar memory modules and/or subsystems) by implementing some or all of the features of the invention described herein in the memory controller (which manages memory module access and proper operation).

[Para 23] [0019] As noted above, the DDR-II does not provide for a 4 bit burst intra-burst interruption, and the minimum unit of data that can be efficiently accessed is 4 bits times the number of I/O pins. Thus in the case of shorter accesses, such as are required for networking operations and other applications, only a fraction of the bursted data is required at the processor level, and the balance represents wasted clock cycles and unproductive latencies. Further, the emerging hyperthreading, multithreading, multichannel-DRAM access, shared memory and other processor inventions referred to above are making multiple open-page capability in every bank an essential capability for the future, without incurring the latencies that currently attend intra-bank, multi-row accesses. In addition, the posted precharge allows the re-accessing of a previous page without incurring worst case scenario latencies.

[Para 24] [0020] In addition to DDR SDRAM's, the Posted Precharge approach of the present invention is equally applicable to FCRAM™, RLDRAM™, RDRAM™, XDRAM™, Network DRAM™, and other RAM architectures utilizing precharge cycles. The multiple open-page architecture of this invention is also

equally applicable to volatile RAM's and nonvolatile RAM's, including flash memory.

[Para 25] [0021] Other objects and advantages of this invention will be better appreciated from the following detailed description.

[Para 26] BRIEF DESCRIPTION OF THE DRAWINGS

[Para 27] [0022] Figure 1 illustrates a block diagram of a commercially-available 64Meg x 4, DDR2 SDRAM.

[Para 28] [0023] Figure 2 illustrates a block diagram of a commercially-available 32Meg x 8, DDR2 SDRAM.

[Para 29] [0024] Figure 3 illustrates a block diagram of a commercially-available 16Meg x 16, DDR2 SDRAM.

[Para 30] [0025] Figure 4 illustrates a block diagram of a commercially-available 8Meg x 32, RL DRAM (reduced latency).

[Para 31] [0026] Figure 5 illustrates a block diagram of a commercially-available 8Meg x 32, common I/O RL DRAM II.

[Para 32] [0027] Figure 6 illustrates a block diagram of a commercially-available 8Meg x 32, separate I/O RL DRAM II.

[Para 33] [0028] Figure 7 illustrates a block diagram of the DRAM device of Figure 1, modified in accordance with an embodiment of the present invention.

[Para 34] [0029] Figure 8 illustrates a block diagram of the DRAM device of Figure 2, modified in accordance with a second embodiment of the present invention.

[Para 35] [0030] Figure 9 illustrates a block diagram of the DRAM device of Figure 3, modified in accordance with a third embodiment of the present invention.

[Para 36] [0031] Figure 10 illustrates a block diagram of a network DRAM modified in accordance with a fourth embodiment of the present invention.

[Para 37] [0032] Figure 11 illustrates an exemplary timing diagram for a standard DDR2 DRAM with burst read and auto-precharge in accordance with the prior art.

[Para 38] [0033] Figure 12 illustrates a timing diagram for a standard DDR2 DRAM modified to have a Posted Precharge in accordance with the present invention.

[Para 39] DETAILED DESCRIPTION OF THE INVENTION

[Para 40] [0034] The utilization of the present invention applies to volatile as well as nonvolatile memories. Implementation in stand alone memory devices, SOC (System On Chip), SIP (System In Package), SIC (System In Chip), DIMM's (Dual In Line Memory Modules), SIMM's (Single In Line Memory

Modules) and other combinations are possible. Furthermore, “page” architecture is widely used in DRAM’s and flash memories, the operations of the latter being described in detail in the available literature and therefore will not be discussed in any detail here. “Precharge” is widely used for dynamic devices like DRAM’s, FeRAM’s (ferroelectric RAM’s), etc. “Page” architecture is also expected to influence future memory products like MAG RAM’s, plastic RAM’s, CNT RAM’s (carbon nano tube), organic memories, phase-change memories, molecular memories and similar products. As such, the implementation of the present invention encompasses all such devices, as well as other memory devices that employ a page architecture.

[Para 41] [0035] Figures 1 through 6 illustrate block diagrams of high level architectures for existing DRAM’s commercially available from Micron Technology, Inc. These block diagrams are merely intended to be representative of known DRAM architectures, and not a limitation to the discussion and application of the present invention. Figure 7 depicts one embodiment of the invention, represented as a modification of the architecture of Figure 1. In Figure 7, a precharge counter (operable to count (system) clock cycles) is shown as being incorporated in the row path of the memory architecture. The counter is similar in its operation to counters employed in VLSI design, though its function and utility are applied to provide a Posted Precharge for the DRAM memory device of Figure 7. More specifically, while counters are used in DDR DRAM’s to refresh data by row and prefetch “burst” bits in the column path, the precharge counter of the present invention is placed in the row path to perform a Posted Precharge operation as described below.

[Para 42] [0036] The precharge counter of this invention has two basic functions. First, when a row address is latched (as a result of a Bank Active Command) and a page is opened, the counter locks into that row address until reset. Second, when a Posted Precharge command is asserted, an internal activation for precharge after ‘n’ number of cycles is activated. The value of n

can be programmed or fixed. Alternatively, n could be set in the Mode Register Set (MRS). Unlike current DRAM's that employ an Auto-Precharge command to automatically close a page, the Posted Precharge of this invention enables an open page (P1) to remain open (available) through the use of latches coupled to the sense amplifier associated with the bank on which the page is located, thereby permitting the storage of data read-from or written-to the sense amplifier. For example, a page can be kept open for 100 cycles or more. In this manner, the precharge counter of this invention enables the 'current page in a specific bank (P1) to be open' for a set time, while permitting the activation of another bank in the same IC and the opening of a different page in a different bank. In addition, by issuing an appropriate command, the current page OPEN time can be extended further (without violating other constraints, such as refresh), by interrupting a Posted Precharge activation (internally) if the memory system decides to extend the current open page.

[Para 43] [0037] In view of the above, until n number of cycles is completed on the precharge counter, if a need arises to go from a current page P2 to a previously opened page (P1), the previously opened page is available as a result of being held open by the precharge counter. After n cycles, if a new row is to be opened in the same bank, the bank goes into precharge after the page (P1) is closed. The present invention further offers the ability to reset the precharge counter if the same row is accessed in a consecutive bank activation cycle. In this manner, the 'n cycle open page' can be extended for as long as the memory system requires it. The precharge (internal) delay provided by the precharge counter – namely, from the time the precharge command is posted to the time the precharge for that particular bank is initiated – can be programably set to any number of desired clock cycles (to maximize bus efficiency).

[Para 44] [0038] Figure 8 shows an embodiment of the invention similar to that of Figure 7, except that an SRAM (static random access memory) is

inserted next to the sense amplifiers. A benefit of this optional feature of the invention is the ability to achieve an ultra low CAS (column address strobe) latency. In computing systems where SDRAM's are used as system memory, there is an overwhelming imbalance between Reads and Writes (Reads far outweigh Writes), and thus one register dedicated only to Reads is preferred. In communication systems where SDRAM's are used for packet buffering, Reads and Writes are balanced; hence, separate SRAM registers for Reads and Writes are recommended. In communications memories where SDRAM's are used as Table Lookup Memory, Reads dominate Writes. Graphic memory, 3D mapping, texture memory, and search engine memories in general belong to the "unbalanced access" class. Although a CPU, NPU, or their associated chip set/cache memory (both on and off-chip) may contain an SRAM, placing the SRAM on the SDRAM chip itself (as represented in Figure 8) provides unparalleled effectiveness in reducing power consumption and bus turnaround times by avoiding "off-DRAM" transactions.

[Para 45] [0039] Because the Posted Precharge function provided by the present invention allows more than one PAGE OPEN per DRAM IC, bus turnaround times are reduced. Figures 9 and 10 represent application of the invention to additional RAMJ devices, with Figure 10 illustrating the application of the invention to a network centric DRAM. A specialty case of DRAM's is 1T SRAM's or Pseudo-SRAM's that use DRAM cores but a non-multiplexed SRAM interface. These specialty RAM's are used mostly in the general field of network memory where random accesses are the dominating type of access. Pseudo-SRAM's are limited foremost by the row cycle time or the RAS pulse width, and usually employ a Read – Auto–Precharge scheme to close the bank as early as possible (i.e., after the output of data to the I/O buffers) in order to speed up such devices by enabling subsequent accesses to different rows. With this operating scheme, a Precharge occurs in the background while a different row is in the process of being opened. However, a problem arises if a subsequent Read request falls into the same row, in which case the request will collide with the ongoing closing of the bank (Precharge) and cause the

device to malfunction or crash. Advantageously, by delaying (posting) a Precharge with the precharge counter of this invention, any subsequent access of the same bank would find it open and would, therefore, be executed without additional row access latencies while concomitantly pushing out the Precharge further. If the subsequent access were to go to a different row, the Precharge of the first row would occur after the Additive Latency of the Posted Precharge. As such, with the flexible internally-timed Posted Precharge capability of this invention, it is possible to eliminate a notable problem encountered by Pseudo-DRAM's, because a subsequent access to the same page simply results in delaying of the Precharge.

[Para 46] [0040] In another embodiment of the invention, instead of an SRAM register, one can use an additional set of sense amplifiers in each bank. Preferably the sense amplifiers are identical, though this is not essential. The ‘page select addresses’ operate on one set of sense amplifiers, while the SDRAM memory bank operates on the other set. There is only one control for both sense amplifier blocks, so that any ambiguity is eliminated. It should be understood that the same concept can be applied to all other SDRAM's, including future SDRAM's comprising more than 4 or 8 banks and SDRAM's of architectures evolving beyond DDR-II.

[Para 47] [0041] Figures 11 and 12 compare the advantages of the present invention to current state-of-the-art DRAM devices, though it will be understood from the foregoing that the “multiple open pages” capability of this invention is easily extendable to flash and other memories. In addition to permitting multiple pages to be held open on different banks, another advantage of this invention is the ability to avoid idle bus cycles. With the memory controller under the supervision of the CPU/NPU (or its chip set), pages can be opened in the memory subsystem SDRAM IC's ahead of processor requirements. Pages can also be closed quickly, such as where a speculative instruction execution does not yield the desired result. A counter

can also be used to keep track of when a page can be closed, so that a Posted Precharge can be Activated for continuous, peak bandwidth operation.

[Para 48] [0042] Another advantage of this invention is the ability to make a previously accessed page available even while a new bank/row address is presented (about 3 clock cycles). This operation improves effective bandwidth when data are written across page boundaries.

[Para 49] [0043] The invention as described above has the ability to solve most known performance issues of SDRAM's (standalone memory devices and memory modules). Each of the disclosed embodiments can be implemented on a memory controller controlling a memory module containing memory devices. It is also possible to design an ASIC to be mounted on such a memory module that contains the functions described above. There are numerous NC (no connect) pins available in commercial SDRAM's for implementing the invention. If it is a 4-bank SDRAM, 2 additional page-select pins can be used to switch among open pages within a SDRAM. An 8-bank SDRAM will require use of 3 pins. Reads/Writes, CAS and other commands require no changes.

[Para 50] [0044] While the invention has been particularly shown and described with reference to particular illustrative embodiments thereof, it will be understood by those skilled in the art that various changes in form and details are within the scope of the invention. Therefore, the scope of the invention is to be limited only by the following claims.

